# Data Mining: Lecture 3

# agenda
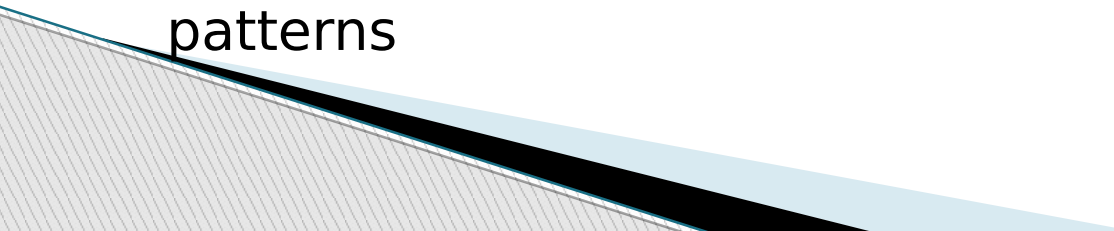
Data Mining Query Language

Data Mining Primitives

# DATA MINING QUERY LANGUAGE (DMQL)

# Why Data Mining Query Language?

- Data mining task may be implemented using a data mining query

- To generate a data mining query, we require data mining task primitives

- Automated vs. query-driven?
  - Finding all the patterns autonomously in a database?—unrealistic because the patterns could be too many but uninteresting

- Data mining should be an interactive process
  - User directs what to be mined

- Users must be provided with a set of primitives to be used to communicate with the data mining system

- Incorporating these primitives in a data mining query language
  - More flexible user interaction
  - Foundation for design of graphical user interface
  - Standardization of data mining industry and practice

# Primitives that Define a Data Mining Task

- Primitive 1: Task-relevant data
  - Database or data warehouse name
  - Database tables or data warehouse cubes
  - Condition for data selection
  - Relevant attributes or dimensions
  - Data grouping criteria
- Primitive 2: Type of knowledge to be mined
  - Characterization, discrimination, association, classification, prediction, clustering, outlier analysis, other data mining tasks
- Primitive 3: Background knowledge
- Primitive 4: Pattern interestingness measurements
- Primitive 5: Visualization/presentation of discovered patterns

# Primitive 3: Background Knowledge

- A typical kind of background knowledge: Concept hierarchies
- Schema hierarchy
  - E.g., street < city < province_or_state < country
- Set-grouping hierarchy
  - E.g., {20-39} = young, {40-59} = middle_aged
- Operation-derived hierarchy
  - email address: hagonzal@cs.uiuc.edu
    - login-name < department < university < country
- Rule-based hierarchy
  - low_profit_margin (X) <= price(X, $P_1$) and cost (X, $P_2$) and ($P_1$ - $P_2$) < $50

# Primitive 4: Pattern Interestingness Measure

- Simplicity
  - e.g., (association) rule length, (decision) tree size
- Certainty
  - e.g., confidence, P(A|B) = #(A and B)/ #(B), classification reliability or accuracy, certainty factor, rule strength, rule quality, discriminating weight, etc.
- Utility
  - potential usefulness, e.g., support (association), noise threshold (description)
- Novelty
  - not previously known, surprising (used to remove redundant rules, e.g., Illinois vs. Champaign rule implication support ratio)

# Primitive 5: Presentation of Discovered Patterns

- Different backgrounds/usages may require different forms of representation

  - E.g.,  rules, tables, crosstabs, pie/bar chart, etc.

- Concept hierarchy is also important

  - Discovered knowledge might be more understandable when represented at high level of abstraction

- Different kinds of knowledge require different representation: association, classification, clustering, etc.

# DMQL—A Data Mining Query Language for Teaching

- Motivation
  - A DMQL can provide the ability to support ad-hoc and interactive data mining
  - By providing a standardized language like SQL
    - Hope to achieve a similar effect like that SQL has on relational database
    - Foundation for system development and evolution
    - Facilitate information exchange, technology transfer, commercialization and wide acceptance
- Design
  - DMQL is designed with the primitives described earlier

# An Example Query in DMQL

**Example 1.11 Mining classification rules.** Suppose, as a marketing manager of *AllElectronics*, you would like to classify customers based on their buying patterns. You are especially interested in those customers whose salary is no less than $40,000, and who have bought more than $1,000 worth of items, each of which is priced at no less than $100. In particular, you are interested in the customer's age, income, the types of items purchased, the purchase location, and where the items were made. You would like to view the resulting classification in the form of rules. This data mining query is expressed in DMQL[3] as follows, where each line of the query has been enumerated to aid in our discussion.

```
use database AllElectronics_db
use hierarchy location_hierarchy for T.branch, age_hierarchy for C.age
mine classification as promising_customers
in relevance to C.age, C.income, I.type, I.place_made, T.branch
from customer C, item I, transaction T
where I.item_ID = T.item_ID and C.cust_ID = T.cust_ID
        and C.income ≥ 40,000 and I.price ≥ 100
group by T.cust_ID
having sum(I.price) ≥ 1,000
display as rules
```

# Other Data Mining Languages & Standardization Efforts

- Association rule language specifications
  - MSQL (Imielinski & Virmani'99)
  - MineRule (Meo Psaila and Ceri'96)
  - Query flocks based on Datalog syntax (Tsur et al'98)
- OLEDB for DM (Microsoft'2000) and recently DMX (XML styled data mining language)
  - Based on OLE, OLE DB, OLE DB for OLAP, C#
  - Integrating DBMS, data warehouse and data mining
- DMML (Data Mining Mark-up Language) by DMG (www.dmg.org)
  - Providing a platform and process structure for effective data mining
  - Emphasizing on deploying data mining technology to solve business problems
- CRISP-DM
  - CRoss Industry Standard Process for Data Mining

# Integration of Data Mining and Data Warehousing

- **Data mining systems, DBMS, Data warehouse systems coupling**

  - No coupling, loose-coupling, semi-tight-coupling, tight-coupling

- **On-line analytical mining data**

  - integration of mining and OLAP technologies

- **Interactive mining multi-level knowledge**

  - Necessity of mining knowledge and patterns at different levels of abstraction.

- **Integration of multiple mining functions**

  - Characterized classification, first clustering and then association

# Coupling Data Mining with DB/DW Systems

- No coupling—flat file processing, no function of DB/DW is utilized, not recommended
- Loose coupling
  - Fetching data from DB/DW
  - No data structure or query optimization as the process is main memory-based
- Semi-tight coupling—enhanced DM performance
  - Provide efficient implementation of a few data mining primitives in a DB/DW system, e.g., sorting, indexing, aggregation, histogram analysis, multiway join, precomputation of some stat functions
- Tight coupling—A uniform information processing environment
  - DM is smoothly integrated into a DB/DW system, mining query is optimized based on mining query, indexing, query processing methods, etc.